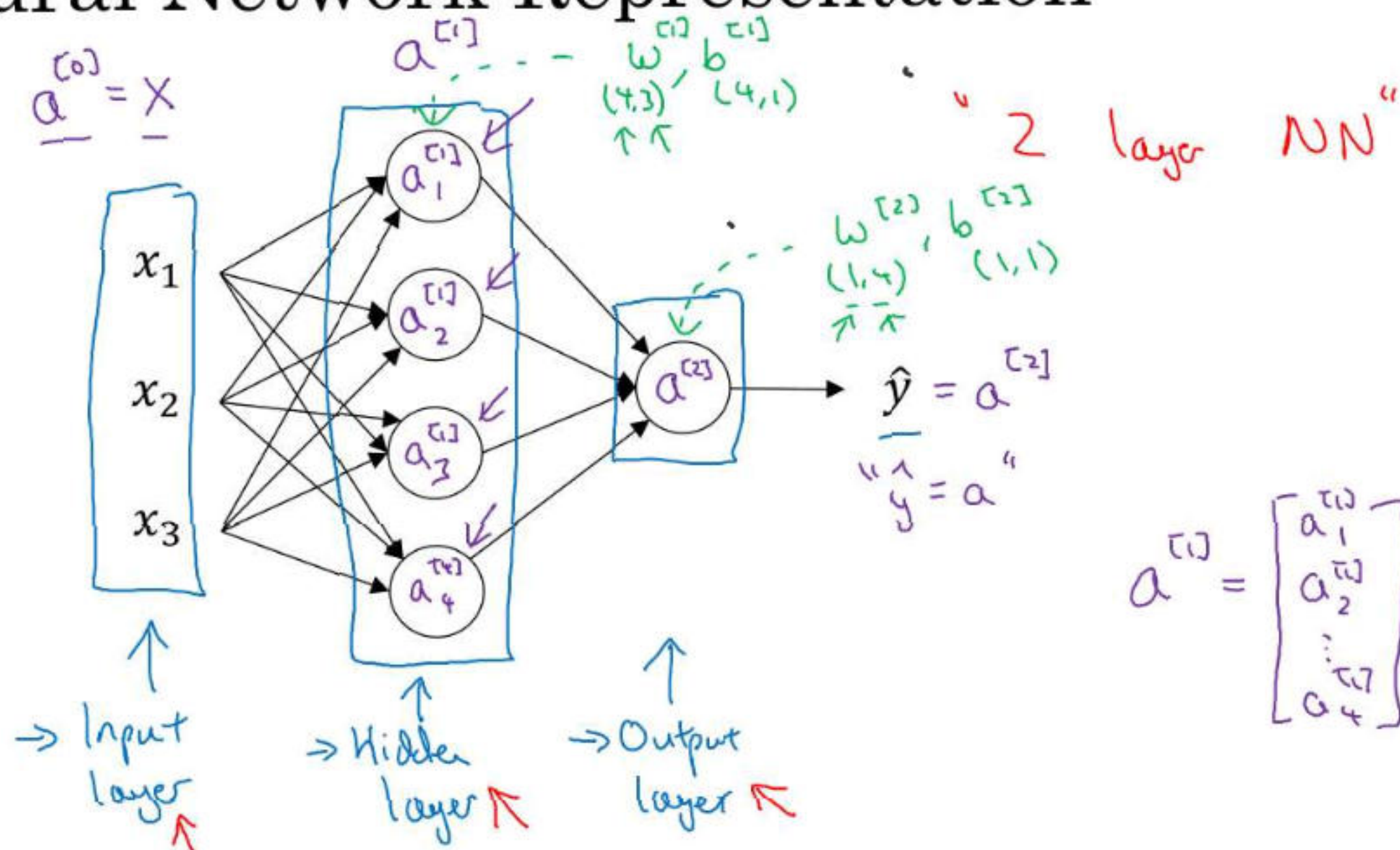


## Neural Network Representation



Andrew Ng

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$$

$$\vdots$$

$$\begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ \vdots \end{bmatrix} x + \begin{bmatrix} b_1^{[1]} \\ \vdots \\ b_n^{[1]} \end{bmatrix}$$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_n^{[1]} \end{bmatrix}$$

$$z^{[1]} = w^{[1]T} x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = w^{[2]T} \cdot a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

12 15.

多个实例:

每个实例  $x^{(i)}$ , 都会获得  $\Rightarrow a^{(2)(i)} \rightarrow \hat{y}$ .

for  $i=1$  to  $m$ :

$$z^{(1)(i)} = W^{(1)} x^{(i)} + b^{(1)}$$

$$a^{(1)(i)} = \sigma(z^{(1)(i)})$$

$$z^{(2)(i)} = W^{(2)} a^{(1)(i)} + b^{(2)}$$

$$a^{(2)(i)} = \sigma(z^{(2)(i)})$$

$$\text{令 } X = \begin{bmatrix} | & & | \\ x^{(1)} & \dots & x^{(m)} \\ | & & | \end{bmatrix}$$

$$A^{(1)} = \begin{bmatrix} | & & | \\ a^{(1)(1)} & \dots & a^{(1)(m)} \\ | & & | \end{bmatrix}$$

node.

$$z^{(1)} = W^{(1)} X + b^{(1)}$$

$$A^{(1)} = \sigma(z^{(1)})$$

$$z^{(2)} = W^{(2)} A^{(1)} + b^{(2)}$$

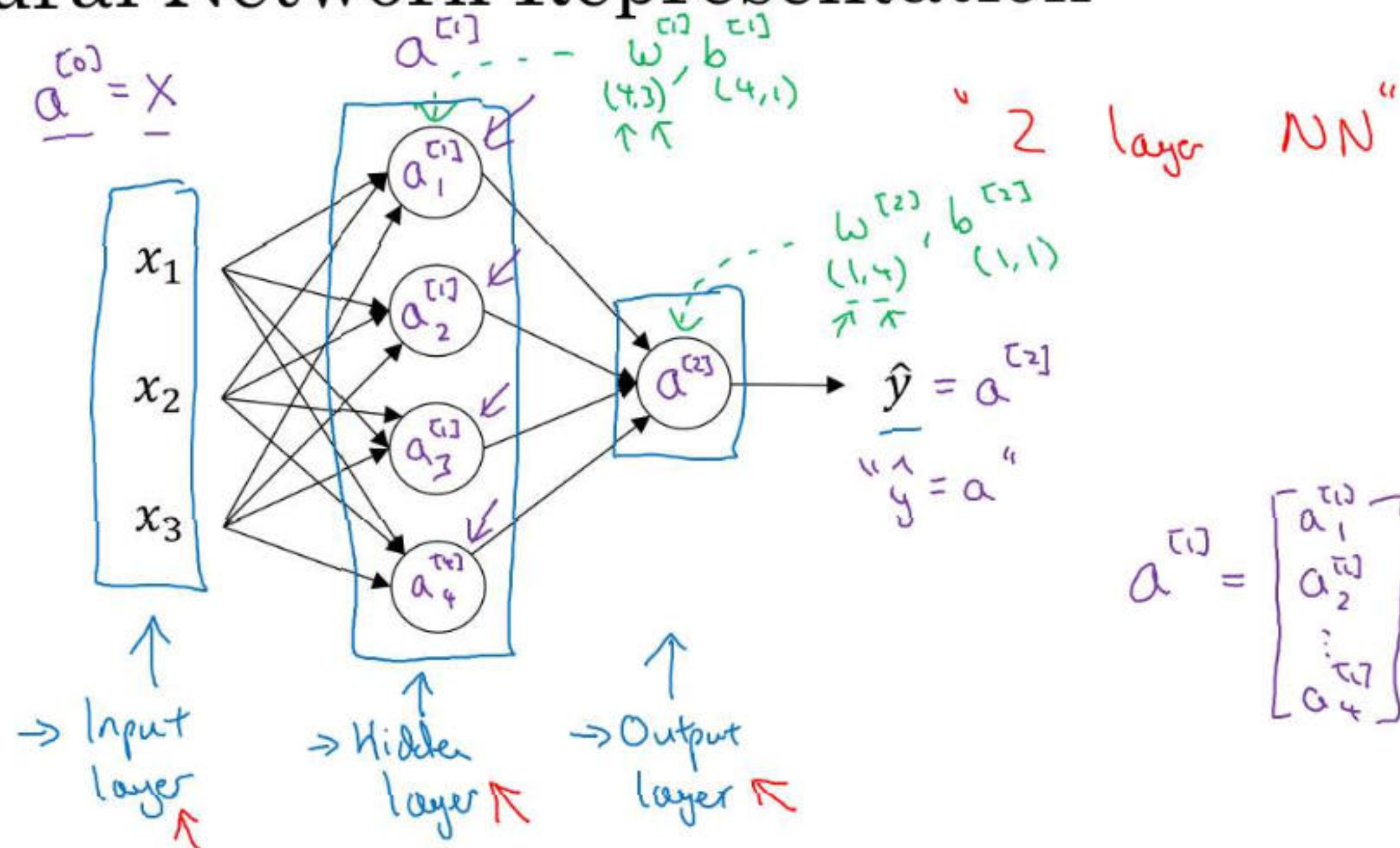
$$A^{(2)} = \sigma(z^{(2)})$$



12 15.

梯度下降:

## Neural Network Representation



Andrew Ng

四个参数:  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$ .

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m L(a^{[2]}_i, y_i)$$

$$\hat{y}^{[1]} = ?$$

$$dw^{[1]} = \frac{\partial J(\theta)}{\partial w^{[1]}}$$

$$w^{[1]} = w^{[1]} - \alpha \cdot dw^{[1]}$$

⋮

$$b^{[2]} = b^{[2]} - \alpha \cdot db^{[2]}$$

$$\left( \frac{1}{1+e^{-z}} \right)' = a(1-a)$$

$$dz^{[2]} = \frac{\partial L}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} = -y(1-a^{[2]}) + (1-y)a^{[2]}$$

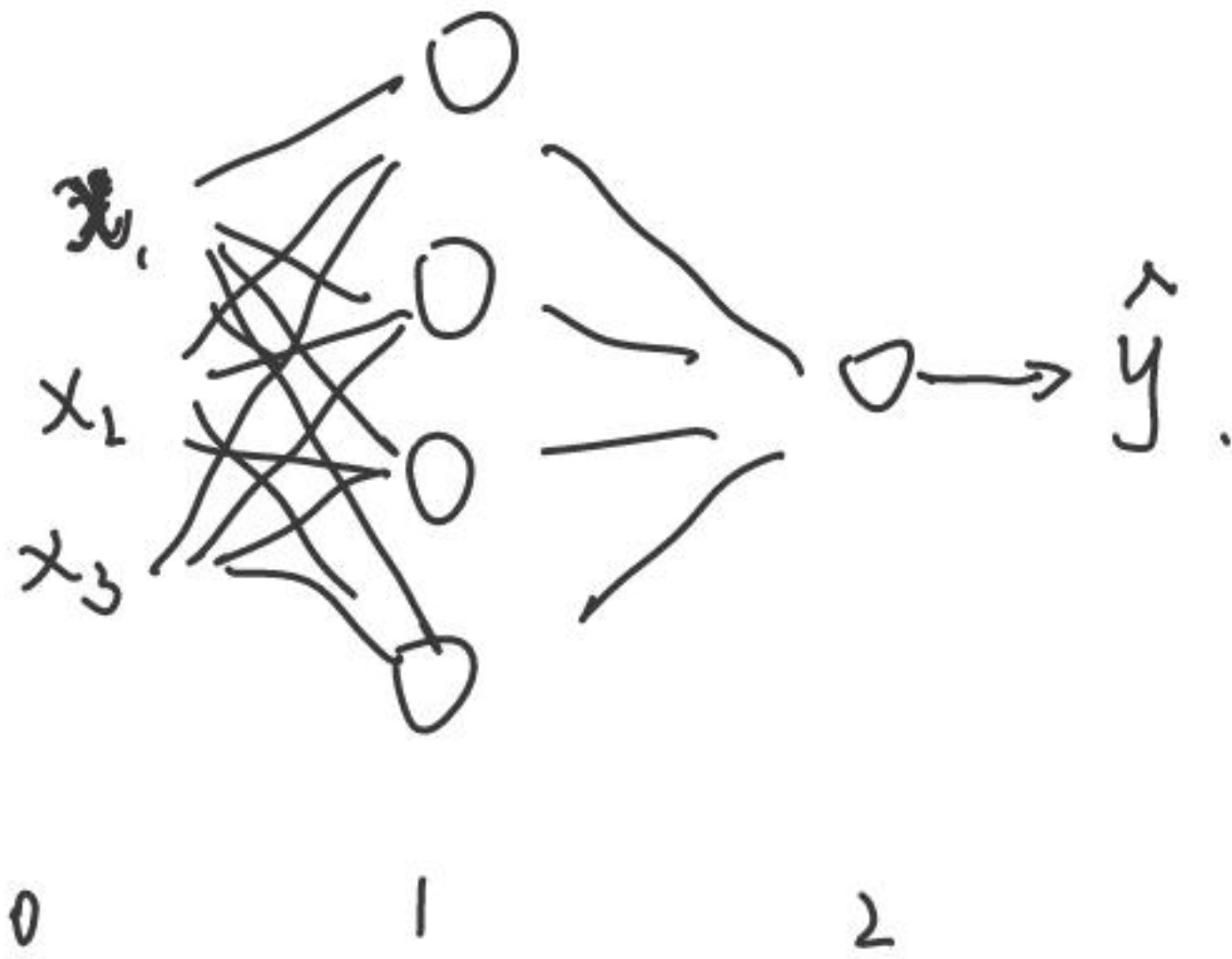
$$dw^{[2]} = dz^{[2]} \cdot \frac{\partial z^{[2]}}{\partial w^{[2]}} = dz^{[2]} \cdot a^{[1]}$$

$$db^{[2]} = dz^{[2]}$$

12 15

$$da^{[1]} = \frac{\partial L}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial a^{[1]}} = dz^{[1]} \cdot w^{[1]}$$

$$dz^{[1]} = da^{[1]} \cdot a^{[1]}(1-a^{[1]})$$



两层有输出

一层 Hidden Layer.



12 16.

## 数据划分.

Train

Dev / Valid

Test.

数据少量: 6:2:2.

大量

Dev + Test < 10%.

Dev, Test 分别有一定数量即可, Test 较少.

## Some strategies.

a) Train 拟合不好:

调整模型, 构建更大 NN, 未收敛? 增加 Iter.

b): Dev 拟合不好:

加 Reg; More Data.

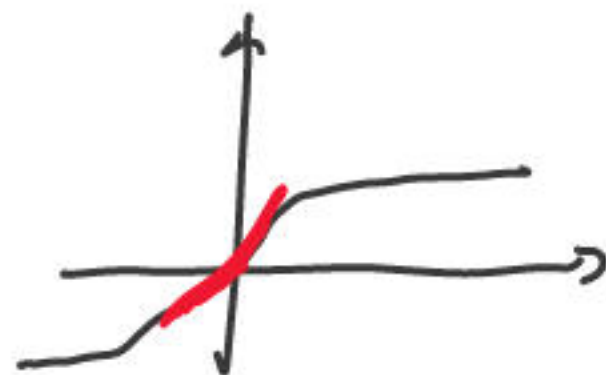
## Regularization:

减少过拟合的两种观点:

a):  $\|w\|_F^2 \approx 0, \Rightarrow$

$z^{[1]} \approx 0,$

$$z^{[1]} = w^{[1]} \cdot a^{[0]} + b^{[1]}.$$



activate:  $g(z)$  近似线性, 线性 NN.

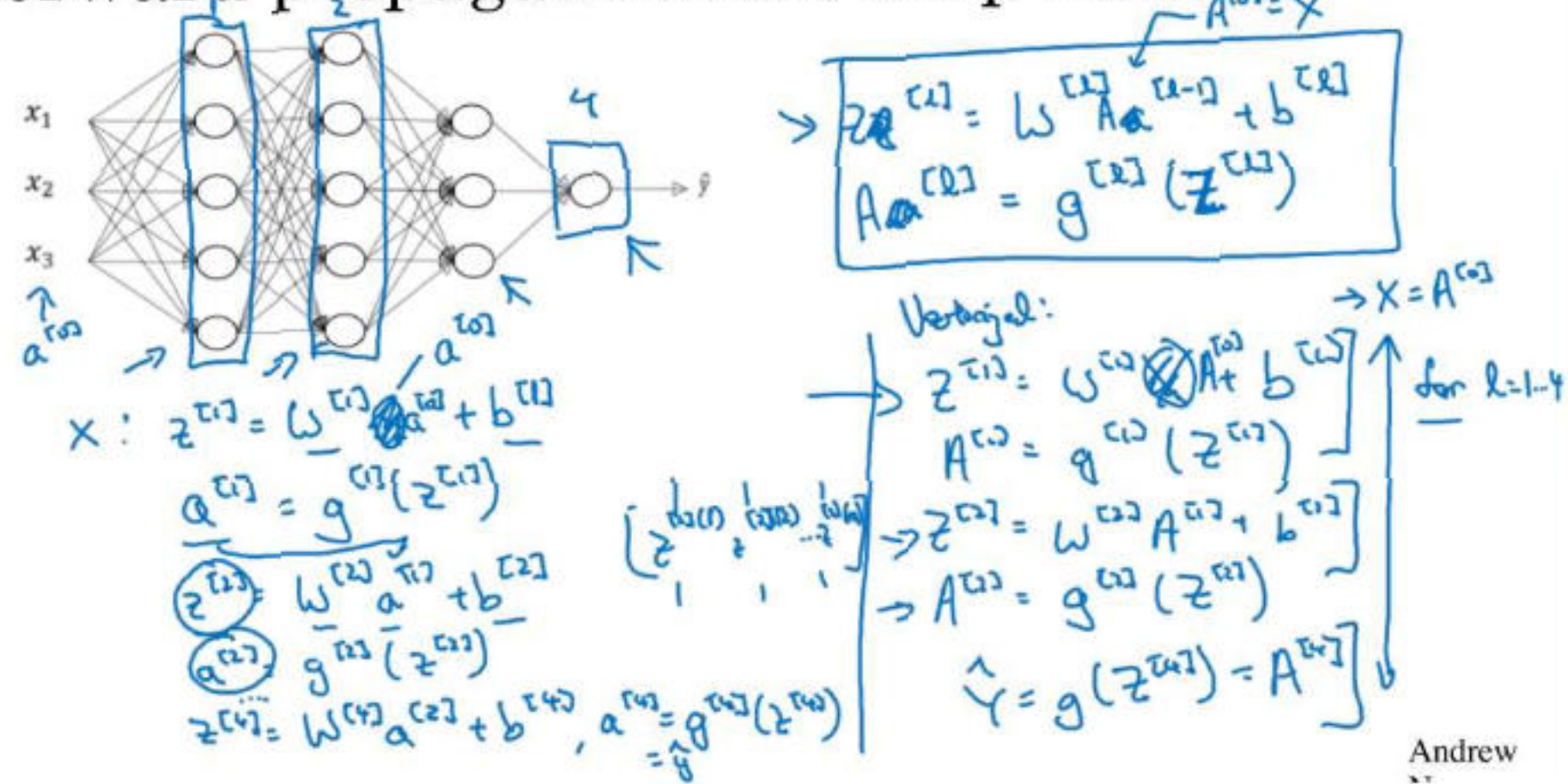
b): 神经元联结近似减少, 简单的 NN.



12 16

# 深层神经网络推导:

## Forward propagation in a deep network



前向传播:

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_s^{[1]} \end{bmatrix} =$$

$$\begin{bmatrix} w_1^{[1]} x + b_1^{[1]} \\ \vdots \\ w_s^{[1]} x + b_s^{[1]} \end{bmatrix} = W^{[1]} x + b^{[1]}$$

$$a^{[1]} = g(z^{[1]})$$

向量化后:

$$z^{[1]} = \begin{bmatrix} z^{[1](1)} & \dots & z^{[1](s)} \\ 1 & & 1 \end{bmatrix} = \begin{bmatrix} w_1^{[1]} x + b_1^{[1]} & \dots & w_s^{[1]} x + b_s^{[1]} \\ 1 & & 1 \end{bmatrix}$$

$$= W^{[1]} \begin{bmatrix} x^{(1)} & \dots & x^{(s)} \\ 1 & & 1 \end{bmatrix} + b^{[1]} = W^{[1]} x + b^{[1]}$$

$$z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[1]} = g(z^{[1]})$$



## Backward propagation for layer $l$

→ Input  $da^{[l]}$

→ Output  $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$\begin{aligned} dz^{[l]} &= da^{[l]} * g^{[l]'}(z^{[l]}) \\ dw^{[l]} &= dz^{[l]} \cdot a^{[l-1]} \\ db^{[l]} &= dz^{[l]} \\ da^{[l-1]} &= W^{[l]T} \cdot dz^{[l]} \\ dz^{[l]} &= W^{[l+1]T} dz^{[l+1]} * g^{[l+1]'}(z^{[l+1]}) \end{aligned}$$

$$\begin{aligned} dz^{[l]} &= dA^{[l]} * g^{[l]'}(z^{[l]}) \\ dw^{[l]} &= \frac{1}{n} dz^{[l]} \cdot A^{[l-1]T} \\ db^{[l]} &= \frac{1}{n} \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{keepdims}=\text{True}) \\ dA^{[l-1]} &= W^{[l]T} \cdot dz^{[l]} \end{aligned}$$

Andrew Ng

反向传播:

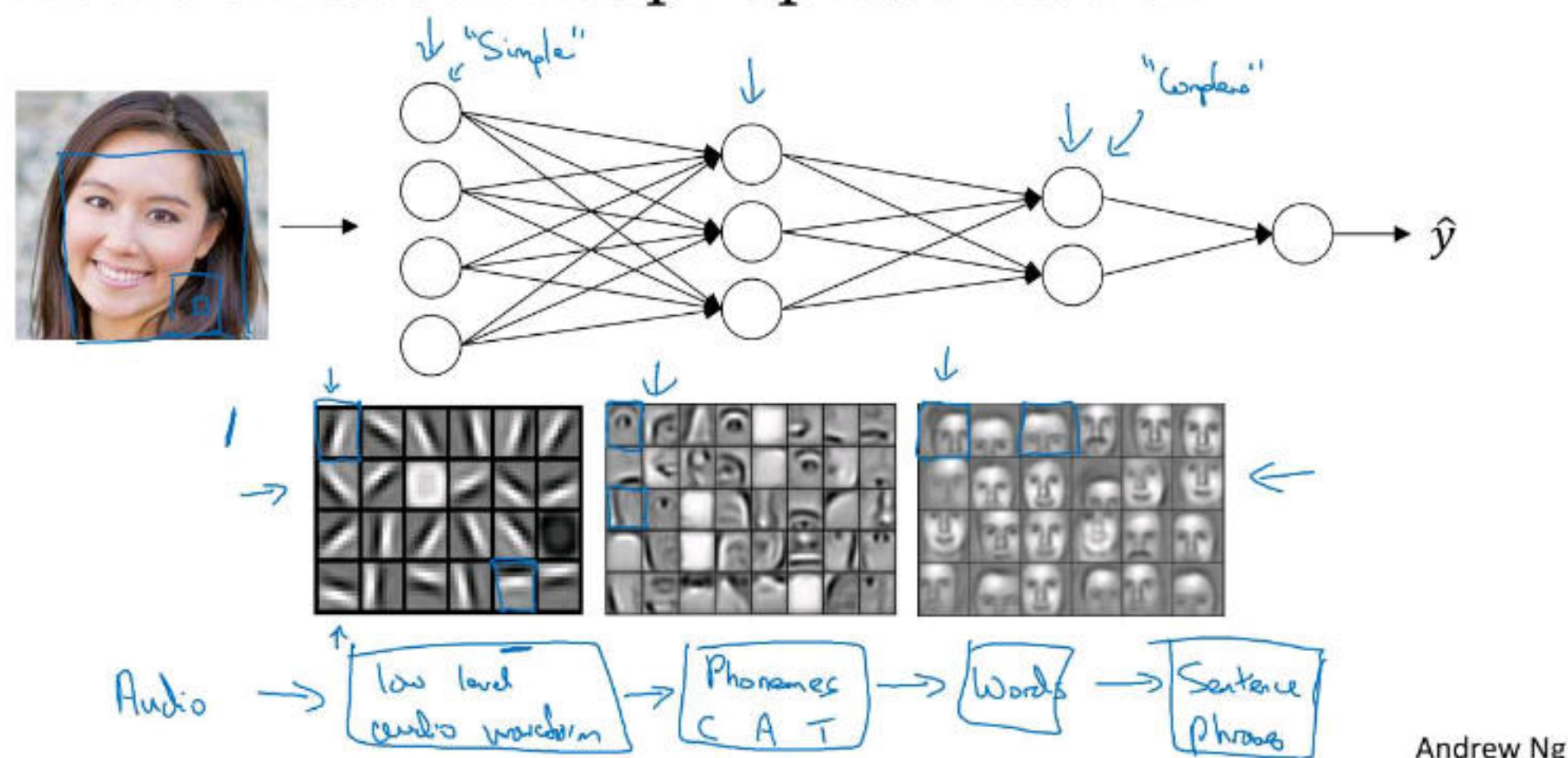
$x_1$	0	0	$a^{[l]}$	
$x_2$	0	0	0	$0 \rightarrow \hat{y}$
$x_3$	0	0	0	
	0	0		

$$\begin{aligned} W^{[l]} \cdot a^{[l-1]} + b^{[l]} &= z^{[l]} & a^{[l]} &= g^{[l]}(z^{[l]}) \\ \frac{\partial J}{\partial z^{[l]}} &= \frac{\partial J}{\partial a^{[l]}} & da^{[l]} &= \frac{\partial J}{\partial z^{[l]}} \cdot \frac{\partial z^{[l]}}{\partial a^{[l-1]}} \\ \frac{\partial J}{\partial z^{[l]}} &= \frac{\partial J}{\partial a^{[l]}} \cdot g^{[l]'}(z^{[l]}) & &= da^{[l]} * g^{[l]'}(z^{[l]}) \\ da^{[l-1]} &= dz^{[l]} \cdot W^{[l]T} & & \end{aligned}$$



12 17

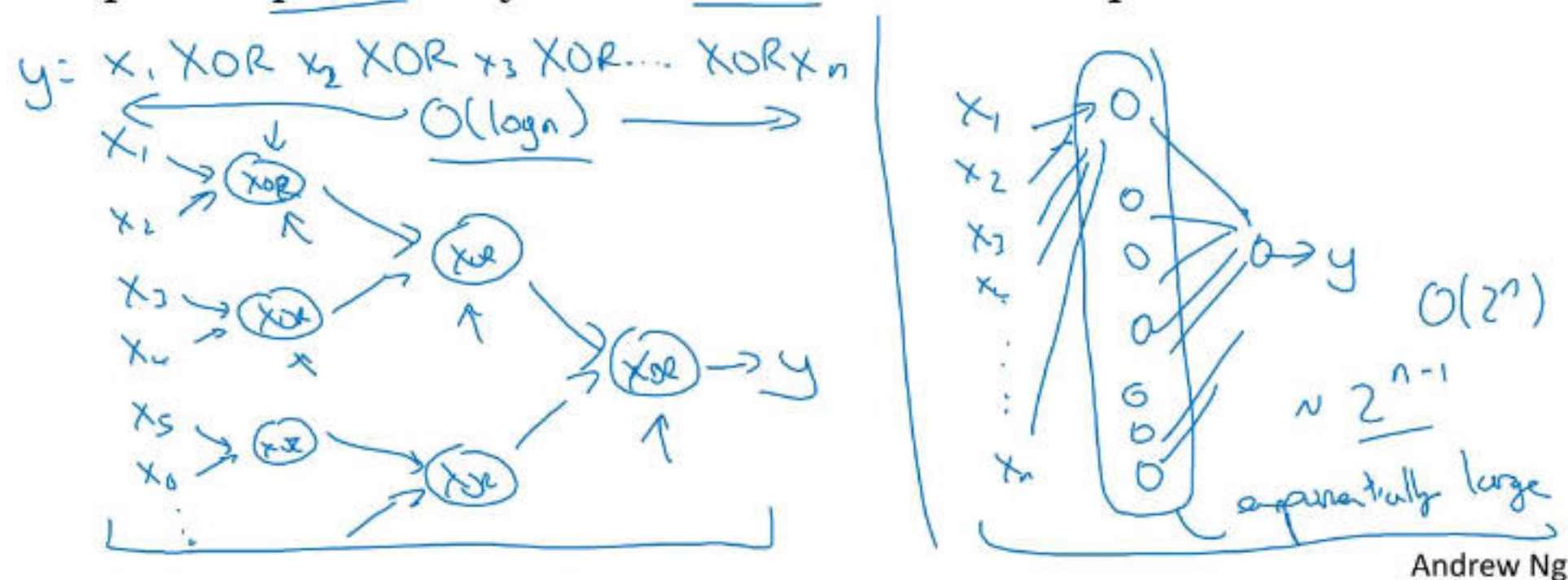
## Intuition about deep representation



Some intuition: 隐藏层提取简单特征, 组合成复杂特征, 最后模式识别.

## Circuit theory and deep learning

Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.



神经元越少 (一个 hidden layer), 需要网络越深. 若只有一层 hidden layer, 神经元的数量会指数级增长.



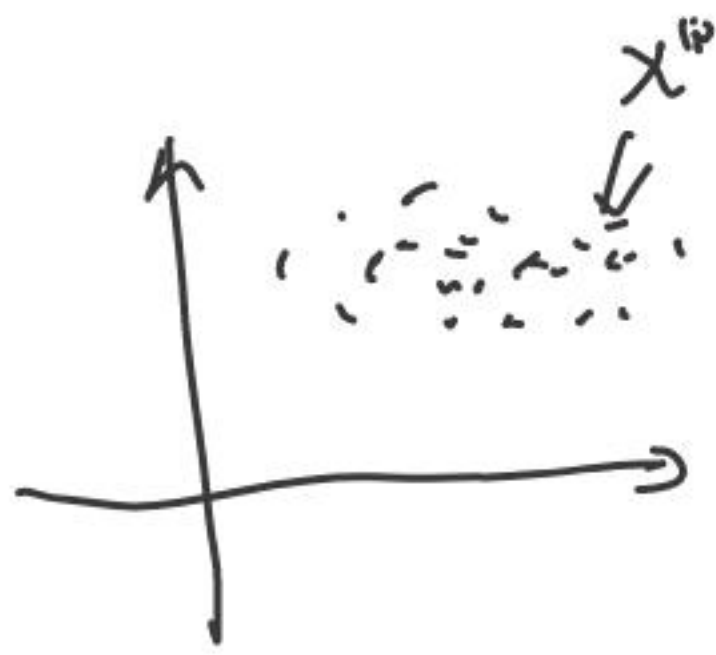
12 17

避免过拟合的其它方法:

a): Data Augmentation: flipped, rotate ...

扩大了训练集.

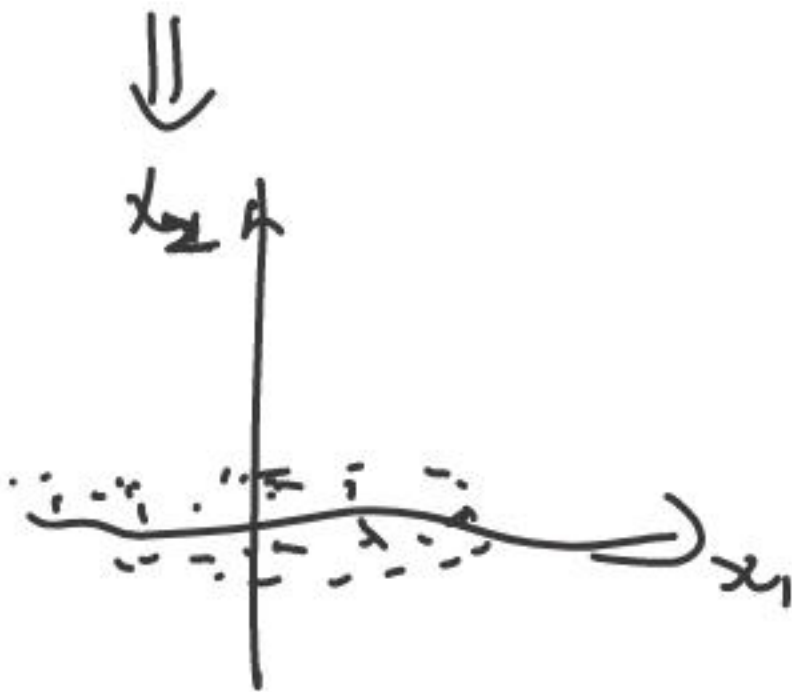
b): Early stopping: 降低 dev set Loss.



归一化特征.

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

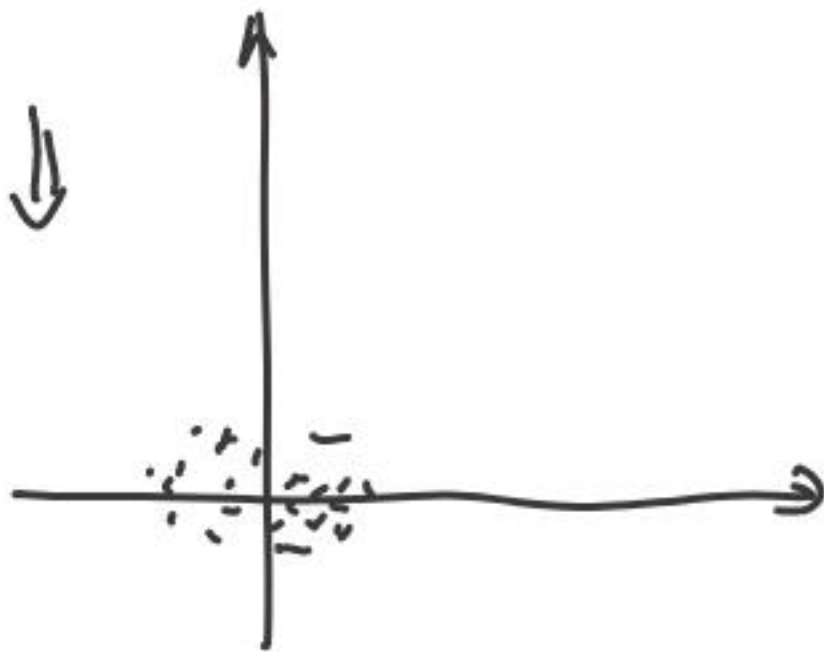
$$x = x - \mu.$$



$$\sigma^2_{x_1} > \sigma^2_{x_2}.$$

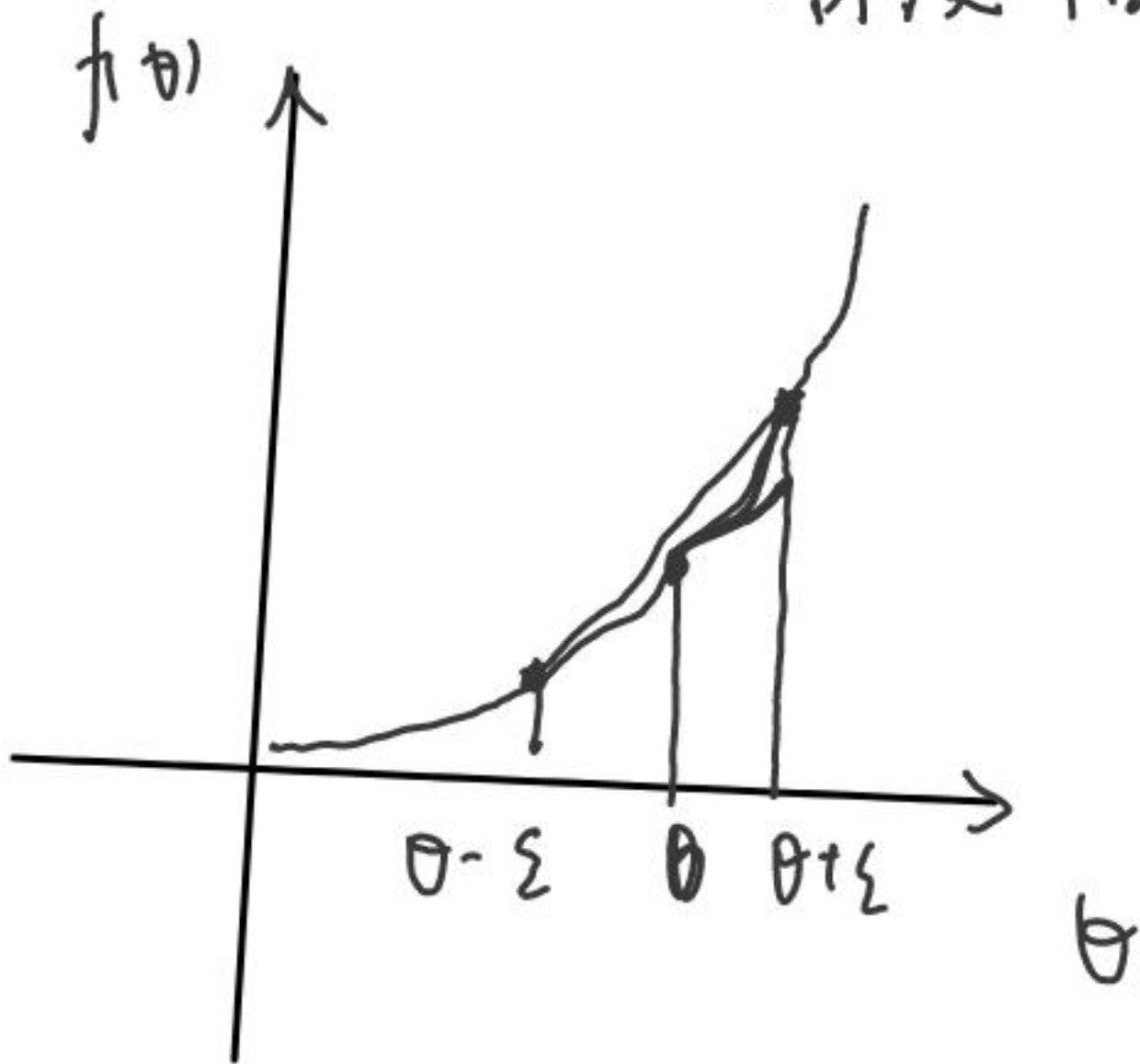
$$\therefore \sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)2}$$

$$x = \frac{x}{\sigma^2}.$$



12 17

梯度估计:



$$f'(\theta) \approx \frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \sim o(\epsilon)$$

$$\approx \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} \sim o(\epsilon).$$

$$f(\theta+\epsilon) = f(\theta) + f'(\theta)\epsilon + \frac{1}{2!}f''(\theta)\epsilon^2 + o(\epsilon^2).$$

$$f(\theta+\epsilon) = f(\theta) + f'(\theta)\epsilon + o(\epsilon).$$

$$f'(\theta) = \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} \sim o(\epsilon)$$



